

AI-driven personalized dance training model and health management research

Tingting Zhang¹ and Hanhua Chen^{1,*}

¹ Institute of Arts, Chongqing College of Humanities, Science & Technology, Chongqing, 401524, China

Corresponding authors: (e-mail: crkart1220099@hotmail.com).

Abstract Under the rapid development of artificial intelligence (AI) technology, this study constructs an AI-driven dance movement personalized training model based on the problems of poor movement recognition accuracy and insufficient personalized instruction in traditional dance training. Methodologically, a two-branch twin supervised learning model is used to realize 2D to 3D skeletal keypoint conversion, and the ST-GCN network is improved by incorporating spatio-temporal attention mechanism to enhance feature extraction in spatial and temporal dimensions. A dataset is constructed using 3,500 images extracted from concert and dance videos, containing six dance movement types, such as crossing the waist, lifting high, spreading one arm, waving, spreading both arms, and walking. The results show that the improved ST-GCN model achieves a recognition accuracy of 93.63% on the test set, which is 14 percentage points higher than the traditional residual network model, and the top-1 metric after fusing spatio-temporal attention is 86.66%, which is 5.63 percentage points higher than the original ST-GCN model. The conclusion shows that the proposed AI-driven dance movement recognition model can effectively solve the problems of movement occlusion and perspective change, significantly improve the recognition accuracy, and provide technical support for personalized dance training and health management.

Index Terms Dance movement recognition, ST-GCN network, Spatio-temporal attention, Personalized training, Health management, Artificial intelligence

I. Introduction

Dance is an art form that requires artists to convey emotions and stories through physical expression [1]. And in dance, the dancer's movement is very important, and the strength, coordination, and coherence it shows directly affects the beauty and charm of the dance [2], [3]. The traditional way of dance movement training is to require dancers to practice through repetition, but dancers have different dance skills, resulting in the overall training effect is not only not obvious, but also a huge time cost [4]-[6]. In addition, in dance training, it is easy to cause dancers to be injured due to training errors, and training with injuries has become the norm [7]. It can be seen that traditional training methods have disadvantages such as slow results and lack of health management, while the application of artificial intelligence (AI) can avoid the above problems from arising [8], [9].

AI is a kind of method and technology by simulating and emulating human intelligence, which has the characteristics of self-learning, self-adaptation, and self-evolution, and can be applied in all walks of life, including dance training [10], [11]. AI can analyze and process the dancer's personalized data to customize the dance training program for each dancer [12], [13]. By collecting and analyzing the dancer's personalized data such as physical fitness, movement mastery, learning style, and dance ability, AI can develop the most suitable dance training plan for the dancer according to his/her actual situation, which can not only effectively improve the effect of dance movement training, but also provide effective health management for the dancer [14]-[17].

This study proposes an AI-driven design scheme for personalized training model of dance movement. First, the accurate conversion of 2D image to 3D skeletal keypoints is realized by constructing a two-branch twin supervised learning network, which effectively solves the problems of movement occlusion and perspective change in traditional methods. Second, the spatio-temporal attention mechanism is integrated on the basis of the classical ST-GCN network to enhance the feature extraction ability in the time and space dimensions, respectively, and improve the recognition accuracy of complex dance movements. Finally, combining the health management characteristics of dance movement, we design corresponding psychological adjustment and physical training modules to construct a complete personalized training system. Through validation on a dataset containing six types of dance movements, it is proved that the proposed method outperforms the existing techniques in terms of recognition accuracy and computational efficiency, and provides an effective technical solution for the construction of intelligent dance teaching and health management system.

II. Dance movement recognition studies

II. A. Recognition algorithm framework

In the dance training or performance scene, the dancer's action gesture changes at a fast speed and with a large amplitude, and there is a certain degree of action masking due to the influence of the performance costume, which further increases the difficulty of action recognition. Traditional action recognition methods utilize optical sensors to collect raw data of action gestures, which are usually video or two-dimensional image data. This type of raw data contains invalid data such as background images and noise interference. Therefore, before action recognition, the raw data need to be pre-processed. After the completion of the original data background removal and other preprocessing, and then extract the skeletal site data, and then realize the data statistics and feature extraction of human posture changes.

Most of the traditional action recognition methods use action data from a single viewpoint or a specific viewpoint, and the model trained with this data has a lower accuracy in recognizing data from other viewpoints. To address the problem of specific dance action recognition, this paper, on the basis of traditional action recognition methods, realizes 3D pose estimation of 2D skeletal keypoints by constructing a 2D skeleton to 3D skeleton two-branch network, thus solving the action occlusion problem. The framework of the specific dance action recognition method proposed in this paper is shown in Fig. 1.

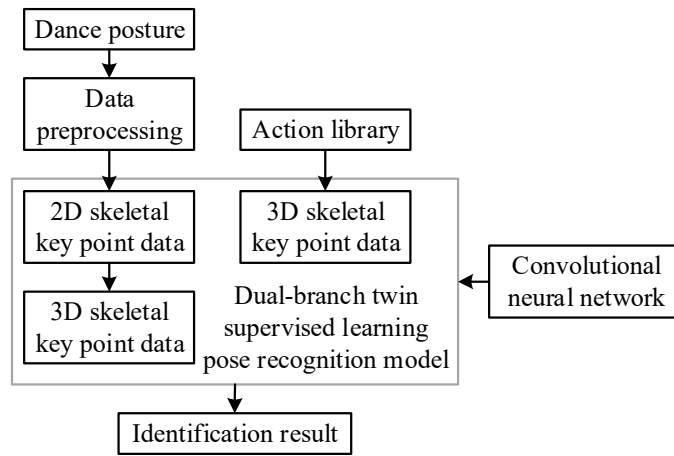


Figure 1: Specific dance action recognition algorithm framework

II. B. Dance Movement Recognition

II. B. 1) Image data preprocessing

At present, the common sensors are optical sensors, the collected image data are color, and the background complexity of the screen is different in different scenes, direct action gesture recognition will increase the computational volume of the algorithm and recognition difficulty. Therefore, it is necessary to carry out pre-processing first, the specific operations are image grayscaling, image thresholding and the construction of human body model refinement. Color images are usually RGB images, which are three-dimensional images in color, and the three-dimensional image can be reduced to one-dimensional by grayscaling, so as to reduce the amount of data. In order to further reduce the amount of data for image processing, it is also necessary to separate the human body contour from the background, and the one-dimensional maximum entropy thresholding method is used in the paper for image background removal. Considering that there are certain differences in the human body contours of different genders, heights and weights, the human body model also needs to be simplified when performing action pose recognition.

The representation of the human body gesture is described using the key point distance matrix M as follows:

$$M = \begin{bmatrix} d_{1,2}^1 & \cdots & d_{k,(k-1)}^1 \\ \cdots & \cdots & \cdots \\ d_{1,2}^N & \cdots & d_{k,(k-1)}^N \end{bmatrix} \quad (1)$$

where, k represents the number of keypoints extracted from the image of this action pose and the maximum value is 14. N represents the number of images of this action pose. $d_{i,j}^1$ represents the distance between keypoints i and j in the 1st image of this action pose.

II. B. 2) Skeletal Posture Critical Point Estimation Models

As the problems of overlapping and occlusion of keypoints inevitably occur in 2D images, and there is a certain similarity between the 2D images projected from different spatial relations of poses in a specific viewpoint. To address this problem, a two-branch twin supervised learning model is designed in the paper to transform 2D images into 3D space for estimating human skeleton keypoints. The 2D key point location pixels of the skeleton in the image are detected by the Open Pose human skeleton 2D pose detector, and the 2D key point pixel matrix P_{2D} is obtained. The two-branch twin supervised learning model, on the other hand, performs the pose estimation of the 3D skeleton key points by building a 2D to 3D key point transformation function [18], which can be described as follows:

$$P_{3D} = X_{change}[P_{2D}, p(c)] \quad (2)$$

$$p(c) = \{p_1, \dots, p_n\} \quad (3)$$

where $p(c)$ is the training parameter of the two-branch supervised learning model, and $X_{change}(\cdot)$ is the loss function of this learning model, which is trained to minimize the gap between the 3D pose estimates and the 3D distribution of the real skeleton keypoints. This time, the distance D is used to characterize the gap between the prediction and the actual 3D skeletal pose position coordinates P_{3DT} as:

$$D_{\min} = \arg \sum_{n=1}^N x_{loss}[X_{change}(P_{2D}), P_{3DT}] \quad (4)$$

where x_{loss} is the loss function for predicting 3D skeletal pose coordinates.

In order to avoid model overfitting, the two-branch twin supervised learning model is improved in the paper: a linear layer, batch normalization layer in convolutional neural network (CNN) is added on each branch. The inputs to the two two-branch networks are preprocessed 2D skeletal keypoints and real 3D skeletal keypoint data from the Human3.6M dataset, respectively. The two branch networks are weighted and fused after the training is completed to finally obtain the 3D skeletal pose keypoint prediction data.

In order to provide accurate recognition of action postures, in addition to recognizing the spatial relationship of skeletal keypoints, temporal coherence is also extremely important. Therefore, the recognition model of 3D skeletal posture needs to be optimized in both spatial and temporal dimensions.

The above model needs to be estimated based on how the skeletal key points are related to each other as follows:

$$q_{at,j} = \sum_{i \in p} w_{i,j} q_{t,i} \quad (5)$$

where, $q_{at,j}$ represents the coordinate information of key point j in the t nd frame image obtained based on spatial dimension. w_i , represents the optimized weight of key point i on the coordinates of j , and $q_{i,j}$ represents the set of all key points connected to the corresponding key point.

In order to enhance the coherence of action gesture recognition at the time level, it is necessary to accurately delineate the start and end frames of the action in a series of images. Moreover, since the images are acquired at certain time intervals, in the case of rapid changes in the action gesture, the position of the 2D skeletal keypoints in the images of the neighboring frames varies greatly, which is prone to misjudgment of the continuity of the action gesture. In the paper, the position change speed relationship between the front and back frames is utilized to enhance the correlation between the front and back frames of the action gesture, as shown in the following equation:

$$q_{mt,j} = \frac{1}{2}(q_{t+1,j} + q_{t-1,j} + \alpha \Delta v_{bt,j} + \beta \Delta v_{ft,j}) \quad (6)$$

$$\Delta v_{bt,j} = q_{t-1,j} - q_{t-2,j} \quad (7)$$

$$\Delta v_{ft,j} = q_{t+1,j} - q_{t+2,j} \quad (8)$$

where Δv_{hi} characterizes the amount of forward motion change of skeletal key point j , and α represents the weight of the current image to maintain the pose change trend of the previous frame. $\Delta v_{jt,i}$ characterizes the amount of backward motion change of skeletal keypoint j , and β represents the weight of the attitude change trend from the previous frame to the current image. q_{mij} is the 3D skeletal pose coordinate data of the key point j obtained based on the time dimension. Therefore, the coordinates of the final skeletal pose keypoints $q_{t,j}$ can be written:

$$q_{t,j} = w_a \cdot q_{at,j} + w_m \cdot q_{mt,j} \quad (9)$$

where w_a, w_m all represent fusion weights.

II. C. Improved ST-GCN model incorporating spatio-temporal attention

II. C. 1) Spatio-Temporal Graph Convolution ST-GCN

In ST-GCN, the human skeleton map is captured by the method above and transformed into joint point coordinate information, and then the corresponding connection relationship is constructed based on the skeletal data. The constructed human joint point map sequence data is input into the spatio-temporal graph convolution network ST-GCN [19], and the graph convolution in both spatial and temporal dimensions is performed to extract more advanced feature maps. Finally, the final results are output by classifying them with fully connected layers and classifiers. The overall flow is shown in Fig. 2. The joints in the skeleton sequence are modeled in two different dimensions, time and space, respectively. Translated in mathematical language, it can be expressed as inputting a sequence of joint points (N, T) of the human body, where N denotes N joint points of the human body and T denotes the length of the input sequence. An undirected graph $G = (V, E)$ is constructed on this basis, with V representing the set of graph nodes, i.e., $V = \{v_\alpha | t = 1, 2, \dots, T, i = 1, 2, \dots, N\}$. E represents the set of edges, which is composed of two parts, E and E , E , which represents the realistic connection of human joint points in the current frame, and E_t represents the time connection of the same joint points in different frames.

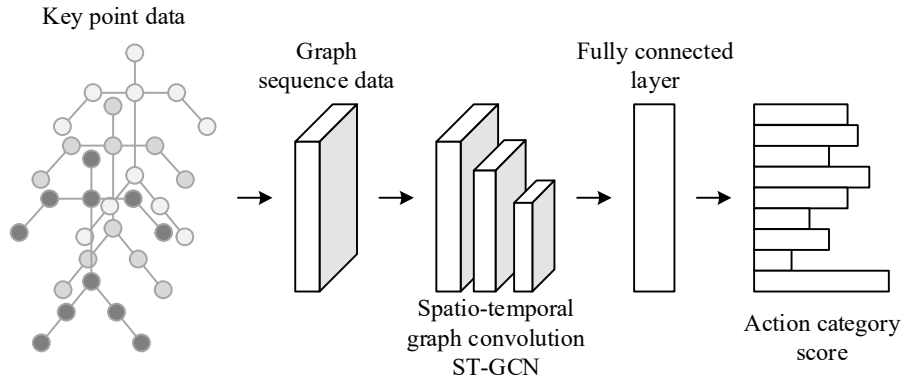


Figure 2: ST-GCN identification process

From the well defined undirected graph G in the above paragraph, the graph convolution operation is then defined under the current frame (at the spatial dimension level). For the node V_{ri} on the τ frame it can be represented as:

$$f_{mat}(v_{ri}) = \sum_{r_n \in s_i} \frac{1}{T_{ij}} f_m(v_{ri}) w(l_i(v_{ri})) \quad (10)$$

where, v represents the nodes of the graph G , f_n represents the feature mapping, s_i is the sampling region for the convolution of the target node v_r , the weight function w is used to provide the weight vector, and the

mapping function l assigns weights to the feature vectors. The size of s_i varies, as does the number of subsets it contains.

Converting the formulae yields the formulae for graph convolution realized in spatial dimensions as:

$$f_{out} = \sum_k^{\kappa_e} w_e (f_m (\tilde{A}_k \odot M_k)) \quad (11)$$

where K denotes the size of the convolution kernel, \tilde{A} is the normalized form of the adjacency matrix A , M is a learnable weight matrix, and the \odot symbol denotes the dot product.

The graph convolution on neighboring frames (i.e., the frame before or after the current frame) is defined next. For each vertex V_{ii} , centered on it, looking one frame forward and one frame backward, there are and only those two identical joints corresponding to it, i.e., in the whole sequence of human joints, each human joint has two fixed neighbor nodes from the time level. Therefore, only two-dimensional convolution of the feature graph output from the model is needed to complete the graph convolution operation in the time dimension.

The original ST-GCN network model contains 10 layers of ST-GCN modules, except for the first ST-GCN module, the last 9 ST-GCN modules include not only graph convolution and time convolution modules, but also residual networks. After the feature map is processed by Polling layer and FC layer, it enters the Softmax classifier and finally outputs the result.

II. C. 2) Improved ST-GCN design incorporating spatio-temporal attention

(1) ST-GCN structure with fused spatio-temporal attention

In this paper, we add the spatio-temporal attention module on the basis of the original model ST-GCN module, which extracts and fuses the human joint point features from both the temporal and spatial dimensions, and strengthens the global feature information of the feature map. The spatio-temporal attention module is up-channeled into the temporal attention sub-module: this sub-module first processes the input through a 1×3 convolutional layer, which is used to learn the weight distribution of the video sequence in the temporal dimension. Then it goes through a Batch Normalization layer for normalizing the data distribution of the output. Finally the output is restricted to a range between 0 and 1 by a Sigmoid activation function, which is used to represent the weights at each time step. Down the road is the spatial attention sub-module, which functions similarly to the temporal attention module, with the difference that a 3×1 convolutional layer is used to process the input, and the final output is the weights for each spatial location. Fusing the two eventually results in a fused feature representation containing both temporal and spatial information. The global feature information of the feature map is enhanced by the learning of the spatio-temporal attention module.

(2) Improved ST-GCN basic unit for fusing spatio-temporal attention

Each basic unit of the ST-GCN model with fused spatio-temporal attention consists of three parts together. That is, it consists of a spatio-temporal attention module, a spatial GCN map convolutional layer, and a temporal TCN convolutional layer. The data first passes through the spatio-temporal attention module, and then enters the spatial map convolution module and the temporal map convolution module to extract the human joint point skeleton feature information after adding temporal and spatial attention in the spatial and temporal dimensions, respectively. Finally, residual links are used to ensure the stability of the training.

The specific calculation process is as follows: the attention weight W_c , σ denotes the sigmoid activation function, Conv2d represents the convolution operation, the size of the convolution kernel is (1, 3), and the padding padding is (0, 1) in the time dimension:

$$W_c = \sigma(\text{Conv2d}(x)) \quad (12)$$

The input x is then subjected to a weighting operation in the time dimension, with \cdot standing for multiplication:

$$x' = x \cdot W_c \quad (13)$$

Similarly, the attention weights W_t are computed in the spatial dimension with a convolution kernel size of (3, 1) and padding padding of (1, 0):

$$W_t = \sigma(\text{Conv2d}(x')) \quad (14)$$

The input x is then weighted on the spatial dimension:

$$x'' = x' \cdot W_t \quad (15)$$

The input feature x'' after the above processing then contains features in both temporal and spatial dimensions.

III. Simulation testing

To verify whether the improved ST-GCN model is effective, training and validation are performed on the public dance movement dataset to evaluate the effectiveness of the algorithm. It is also compared with traditional dance movement recognition methods to verify the recognition efficiency of the model.

III. A. Sources of data sets

The experimental dataset comes from a total of 3500 image frames extracted from concert videos and dance videos. The content of the dataset includes singers and dancers at home and abroad, and the stage includes various kinds of large and small performance platforms and daytime scenes and nighttime scenes. Firstly, the key points of the images are obtained by pose recognition, and then the work features and image features are calculated to get the single-frame images and 18-dimensional dataset of the characters' movements. It contains six kinds of dance movements, such as crossing the waist, raising high, spreading one arm, waving, spreading both arms, and walking. 3200 images in the dataset are randomly selected as the training set and the remaining 300 images are used as the test set.

III. B. Evaluation indicators

The action recognition under study is a classification task, and the main consideration is the correctness of the classification, i.e., the accuracy of the recognition, which is calculated as shown in equation (16). The accuracy rate of the classification index is defined as the ratio of the number of correctly recognized samples to the total number of samples participating in the test. Where TP and TN denote the number of correctly classified samples, and FP and FN denote the number of incorrectly classified samples:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (16)$$

In the field of action recognition research on large-scale datasets, top-1 and top-5 are commonly used as performance metrics for algorithms. Where top-1 indicates the ratio of the number of samples classified correctly in the sample with the highest prediction probability to the number of all samples participating in the prediction, and top-5 indicates the ratio of the number of samples classified correctly in the sample with the top 5 prediction probability to the number of all samples participating in the prediction, and in the experiments in this paper, we take the best result of multiple measurements as the value of the top-1 metric.

III. C. Ablation experiments

To verify the effectiveness of the algorithms proposed in this chapter, this section sets up comparative experiments, which follow the order of modules first, then the whole. The ST-GCN is used as the baseline model for the experiment, and the ablation experiment is carried out on the dataset. Firstly, the effectiveness of temporal attention proposed in this chapter is verified, secondly, the effectiveness of spatial attention is verified, and exploratory experiments and visualization of results on the actual effect are conducted. Finally, the recognition performance of the overall model incorporating temporal and spatial attention is verified on the dataset.

(1) Experiments on the effectiveness of temporal attention

The study uses the top-1 evaluation index for assessment, and the results of temporal attention effectiveness ablation experiments are shown in Table 1. The top-1 reaches 82.63% after adding temporal attention, which is 1.6 percentage points higher than the original.

Table 1: Time attention effectiveness ablation experiment results

Method	Top-1
ST-GCN	81.03%
ST-GCN + Time attention	82.63%

In order to visualize the effect of temporal attention on the model, the connection strengths of the nodes of the initial model and the model after adding temporal attention are demonstrated, and the results of the visual analysis are shown in Fig. 3. (a) and (b) represent the training process of the ST-GCN model and the ST-GCN model after adding temporal attention, respectively. In the dataset, there are 11 nodes in each skeleton, and the color depth of each point the table is in the position of the connection strength of the two nodes. The results in the figure show that

the ST-GCN model designed in this paper after adding temporal attention is able to update the connection weights between the nodes in the adjacency matrix during the training process according to the graph structure of the data samples. At the same time, new generative weights are also given between nodes that are not physically connected spatially to establish virtual connections for two non-neighboring two nodes, which enhances the flexibility of the model for different action samples and thus improves the action recognition accuracy.

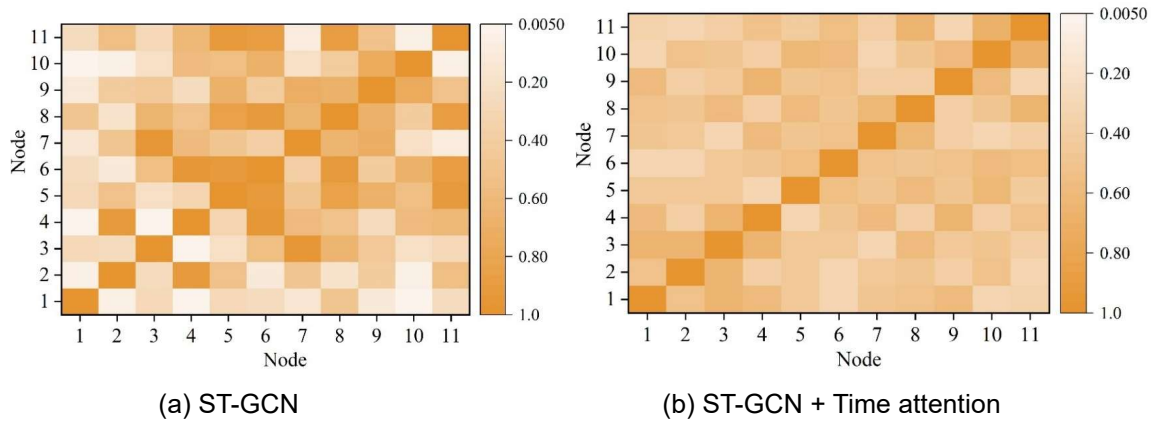


Figure 3: The connection strength of the node

(2) Spatial attention validity experiment

The study was assessed using the top-1 evaluation index, and the results of the spatial attention effectiveness ablation experiment are shown in Table 2. The top-1 reached 82.82% after adding spatial attention, which is 1.79 percentage points higher than the original.

Table 2: The results of the experiment of spatial attention effectiveness

Method	Top-1
ST-GCN	81.03%
ST-GCN + Space attention	82.82%

The connection strength of the nodes of the model after adding spatial attention is demonstrated, and the visualization and analysis results are shown in Fig. 4. The visualization results of ST-GCN model training after adding spatial attention are compared with Fig. 3(a), the color depth is deepened and the connection weights between nodes in the adjacency matrix are strengthened. The flexibility of the model for different action samples is enhanced after adding spatial attention.

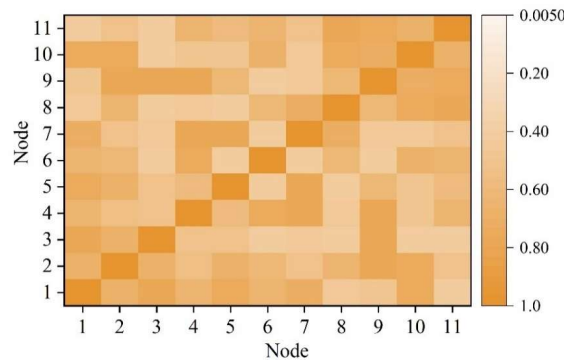


Figure 4: The strength of the connection of ST-GCN + Space attention

(3) Experiment on the effectiveness of fused spatio-temporal attention

In order to verify the effectiveness of the ST-GCN model after fusing spatio-temporal attention, the improved model is compared with the ST-GCN model on the dataset, and the comparison results are shown in Table 3. The top-1 reaches 86.66% after fusing spatio-temporal attention, which is 5.63 percentage points higher than the original.

Table 3: Model comparison results

Method	Top-1
ST-GCN	81.03%
ST-GCN + Blend time and space	86.66%

The confusion matrix for the six dance movements in the above dataset was constructed, and the construction results are shown in Fig. 5. To ensure that the number of each movement dataset is basically the same, then the ST-GCN model incorporating spatio-temporal attention is used for recognition. According to the recognition results, it can be seen that the classification accuracy of all six dance movements reaches more than 90%, and the improved ST-GCN model's can realize the recognition of different dance movements.

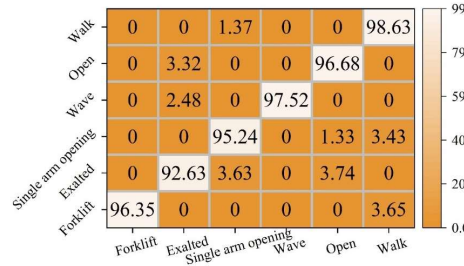


Figure 5: The confusion matrix of six dance movements

In summary, the introduction of spatio-temporal attention to the ST-GCN model allows the network to be more expressive in extracting time-domain features and space-domain features, and to recognize dance movements with higher accuracy.

III. D. Model comparison experiments

In order to verify the superiority of this paper's model, this paper's model is used to perform recognition test with the traditional recognition model residual network four-channel model [20] (Model 1) and computational Hu moments model [21] (Model 2) on the test set, and the comparison results are shown in Table 4. Compared with the comparison model, the accuracy of this paper's model is higher, reaching 93.63%, which indicates that this paper's model is more effective in recognizing dance movements with high accuracy.

Table 4: Contrast model recognition accuracy

Model	Accuracy
Model 1	79.63%
Model 2	81.63%
Ours	93.63%

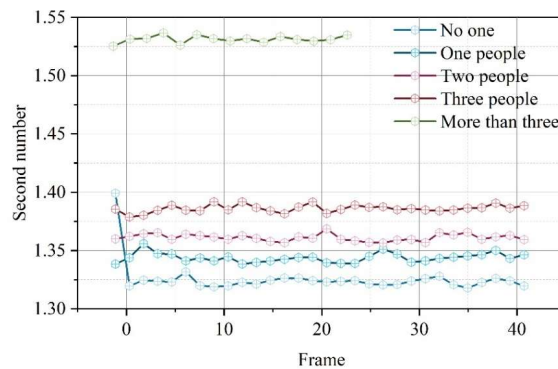


Figure 6: Model efficiency validation

In addition, as shown by the experimental time, the model in this paper runs at 0.76 frames/s on a Tesla P4 graphics card and can recognize multiple people's actions in a single image. In order to further verify the recognition efficiency of this paper's model, it was tested in the scenarios of 0 to multiple people respectively, and it was found that the time spent by this paper's model increased gradually with the increase of the number of people in the image, but it was only a small increase. The running time of this paper's model increases linearly with the increase in the number of people, and the validation of the model efficiency is shown in Fig. 6. In comparison, the running time of the model in this study basically did not increase more substantially. This indicates that the model in this paper is more efficient and the model performance is better.

IV. Health management design

IV. A. Reducing stress and strengthening the psyche

The role of psychotherapy is to create a relaxing and comfortable atmosphere through relaxing and joyful music, collective shouting of slogans and counting of beats, which drives the individual's emotional experience, and at the same time, accompanied by heartfelt body dancing, it can effectively release the pressure in people's hearts and enhance the sense of joyful experience. It can be seen that dance has a very important role in promoting the development of mental health, dance teaching should also reflect the humanization, caring, set up to reduce students' psychological pressure, promote students' mental health teaching content.

IV. B. Peer Interaction and Experiencing Psychology

Socially speaking, Aristotle once said: "Man is a social animal". People rely on the social environment to survive, there is no completely independent individual, and social life is the interaction between people, want to get a foothold in society, also must have interpersonal relationships. In terms of education, constructivism holds that the acquisition of knowledge is not accomplished by transmission, and that knowledge can only be exchanged in an integrated learning situation. In cooperation and communication, the role of "learning community" is fully utilized. The process of peer communication is also an effective means to promote individual cooperation and collective consciousness. Therefore, the design of peer communication in dance teaching content is essential, is a necessary condition for students to carry out social life, and students can experience different interpersonal relationships from interpersonal communication, stimulate the occurrence of emotional experience.

IV. C. Development of thinking and intellectual stimulation

Thinking can be defined as the process of processing and handling information by human beings or animals, including perception, memory, reasoning, judgment, decision-making and many other aspects. Human beings use thinking to solve problems, improve their level of understanding, control their emotions and behaviors and other aspects, it is one of the sources of human wisdom. And thinking ability is the basis of human contact, cognition and problem solving, and the core of human excellent quality and creativity. With the rapid development of information technology and artificial intelligence, many simple and repetitive labor has been replaced by machines. In contrast, talents with more independent thinking ability, creativity and innovative consciousness are more valued by the society and the market, and only with more perfect thinking ability can they understand and evaluate the world more effectively and better cope with the challenges and opportunities in life. Therefore, one of the tasks of school education is to cultivate students' thinking ability, so that they can better cope with various problems and situations in their study and life. It can be seen that the development of thinking ability is one of the most important criteria for measuring the psychological development and psychological maturity of an individual. Therefore, thinking training must also be included in the content of the dance curriculum in order to promote students' intellectual development and psychological maturity.

V. Conclusion

In this study, by constructing an improved ST-GCN network incorporating the spatio-temporal attention mechanism, we successfully realized high-precision recognition and classification of multiple dance movements. The experimental results show that:

The improved model achieves significant performance improvement on a dataset containing 3500 images, and the final recognition accuracy reaches 93.63%, which is 14 percentage points and 12 percentage points higher than the traditional residual network model and Hu moment model, respectively. The two-branch twin supervised learning model effectively solves the key technical difficulties in the conversion of 2D to 3D skeletal keypoints, and significantly improves the understanding of complex dance movements by establishing feature associations in both spatial and temporal dimensions. The introduction of the spatio-temporal attention mechanism enables the network

to adaptively focus on the key action features, and the top-1 index after fusing spatio-temporal attention reaches 86.66%, which proves the effectiveness of the attention mechanism in dance action recognition.

The model shows good computational efficiency in practical applications, with a running rate of 0.76 frames per second on a Tesla P4 graphics card, and supports simultaneous action recognition by multiple people, with a linear growth in running time, which meets the requirements of real-time applications. The constructed confusion matrix of six dance movements shows that the classification accuracy of each movement type is more than 90%, which verifies the stability and reliability of the model. The AI-driven dance movement recognition technology proposed in this study not only solves the limitations of traditional methods in movement masking and perspective change, but also provides core technical support for the construction of personalized dance training and health management system. Combined with the health benefits of dance movement in psychological adjustment, social interaction and intellectual development, this technical solution has a broad application prospect and important social value.

References

- [1] Saumaa, H. (2022). Dance emotions. *Integrative and Complementary Therapies*, 28(3), 134-137.
- [2] Torrents Martín, C., Ric, Á., & Hristovski, R. (2015). Creativity and emergence of specific dance movements using instructional constraints. *Psychology of Aesthetics, Creativity, and the Arts*, 9(1), 65.
- [3] Aristidou, A., Zeng, Q., Stavarakis, E., Yin, K., Cohen-Or, D., Chrysanthou, Y., & Chen, B. (2017, July). Emotion control of unstructured dance movements. In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation* (pp. 1-10).
- [4] Esmail, A., Vranceanu, T., Lussier, M., Predovan, D., Berryman, N., Houle, J., ... & Bherer, L. (2020). Effects of Dance/Movement Training vs. Aerobic Exercise Training on cognition, physical fitness and quality of life in older adults: A randomized controlled trial. *Journal of bodywork and movement therapies*, 24(1), 212-220.
- [5] Burzynska, A. Z., Finc, K., Taylor, B. K., Knecht, A. M., & Kramer, A. F. (2017). The dancing brain: Structural and functional signatures of expert dance training. *Frontiers in human neuroscience*, 11, 566.
- [6] Quinn, M. J., Miltenberger, R. G., & Fogel, V. A. (2015). Using TAGteach to improve the proficiency of dance movements. *Journal of Applied Behavior Analysis*, 48(1), 11-24.
- [7] Fuller, M., Moyle, G. M., Hunt, A. P., & Minnett, G. M. (2019). Ballet and contemporary dance injuries when transitioning to full-time training or professional level dance: a systematic review. *Journal of Dance Medicine & Science*, 23(3), 112-125.
- [8] Long, K. L., Milidonis, M. K., Wildermuth, V. L., Kruse, A. N., & Parham, U. T. (2021). The impact of dance-specific neuromuscular conditioning and injury prevention training on motor control, stability, balance, function and injury in professional ballet dancers: a mixed-methods quasi-experimental study. *International journal of sports physical therapy*, 16(2), 404.
- [9] Dang, Y., Chen, R., Koutedakis, Y., & Wyon, M. A. (2023). The efficacy of physical fitness training on dance injury: a systematic review. *International journal of sports medicine*, 44(02), 108-116.
- [10] Wang, J., & Hao, J. (2025). Analysis of mechanical damage in dance training under artificial intelligence behavior constraints. *International Journal of High Speed Electronics and Systems*, 34(02), 2440086.
- [11] Lin, Z. (2022, May). Research on the risk analysis method of dance technology based on artificial intelligence. In *International Conference on Electronic Information Engineering, Big Data, and Computer Technology (EIBDCT 2022)* (Vol. 12256, pp. 308-313). SPIE.
- [12] Zhang, Y. (2023). Analysis of Transformation of Dance Teaching Model Equipped with AI. *Technology. Frontiers in Educational Research*, 6(31).
- [13] Muangmoon, O. O., Sureephong, P., & Tabia, K. (2017). Dance training tool using kinect-based skeleton tracking and evaluating dancer's performance. In *Advances in Artificial Intelligence: From Theory to Practice: 30th International Conference on Industrial Engineering and Other Applications of Applied Intelligent Systems, IEA/AIE 2017, Arras, France, June 27-30, 2017, Proceedings, Part II 30* (pp. 27-32). Springer International Publishing.
- [14] Wu, L. F. (2023, April). Evaluation on Multimodal Dance Action Recognition Based on Artificial Intelligence Image Technology. In *Proceedings of the 2023 International Conference on Frontiers of Artificial Intelligence and Machine Learning* (pp. 112-117).
- [15] Muneesawang, P., Khan, N. M., Kyan, M., Elder, R. B., Dong, N., Sun, G., ... & Guan, L. (2015). A machine intelligence approach to virtual ballet training. *IEEE MultiMedia*, 22(4), 80-92.
- [16] Wang, Y., & Zheng, G. (2020). Application of artificial intelligence in college dance teaching and its performance analysis. *International Journal of Emerging Technologies in Learning (iJET)*, 15(16), 178-190.
- [17] Xie, Y., Yan, Y., & Li, Y. (2025). The use of artificial intelligence-based Siamese neural network in personalized guidance for sports dance teaching. *Scientific Reports*, 15(1), 12112.
- [18] Yuxiang Yang, Yifan Deng, Jiazhou Li, Meiqi Liu, Yao Yao, Zhaoyuan Peng... & Yingqi Peng. (2024). An Effective Yak Behavior Classification Model with Improved YOLO-Pose Network Using Yak Skeleton Key Points Images. *Agriculture*, 14(10), 1796-1796.
- [19] Qi Lu. (2024). Sports-ACtrans Net: research on multimodal robotic sports action recognition driven via ST-GCN. *Frontiers in Neurobotics*, 18, 1443432-1443432.
- [20] Chao Hao, Cao Yiming & Liu Yongli. (2023). Multi-channel EEG emotion recognition through residual graph attention neural network. *Frontiers in Neuroscience*, 17, 1135850-1135850.
- [21] Zhuang Wu, Shanshan Jiang, Xiaolei Zhou, Yuanyuan Wang, Yuanyuan Zuo, Zhewei Wu... & Qi Liu. (2020). Application of image retrieval based on convolutional neural networks and Hu invariant moment algorithm in computer telecommunications. *Computer Communications*, 150, 729-738.